

Using a Non-Commutative Bernstein Bound to Approximate Some Matrix Algorithms in the Spectral Norm

Malik Magdon-Ismail
 CS Department, Rensselaer Polytechnic Institute,
 Troy, NY 12180, USA.
 magdon@cs.rpi.edu

January 18, 2013

Abstract

We focus on *row sampling* based approximations for matrix algorithms, in particular matrix multiplication, sparse matrix reconstruction, and ℓ_2 regression. For $A \in \mathbb{R}^{m \times d}$ (m points in $d \ll m$ dimensions), and appropriate row-sampling probabilities, which typically depend on the norms of the rows of the $m \times d$ left singular matrix of A (the *leverage scores*), we give row-sampling algorithms with linear (up to polylog factors) dependence on the stable rank of A . This result is achieved through the application of non-commutative Bernstein bounds.

Keywords: row-sampling; matrix multiplication; matrix reconstruction; estimating spectral norm; linear regression; randomized

1 Introduction

Matrix algorithms (eg. matrix multiplication, SVD, ℓ_2 regression) are of widespread use in many application areas: data mining (Azar *et al.*, 2001); recommendations systems (Drineas *et al.*, 2002); information retrieval (Berry *et al.*, 1995; Papadimitriou *et al.*, 2000); web search (Kleinberg, 1999; Achlioptas *et al.*, 2001); clustering (Drineas *et al.*, 2004; McSherry, 2001); mixture modeling (Kannan *et al.*, 2008; Achlioptas and McSherry, 2005); etc. Based on the importance of matrix algorithms, there has been considerable research energy expended on breaking the $O(md^2)$ bound required by exact SVD methods (Golub and van Loan, 1996).

Starting with a seminal result of Frieze *et al.* (1998), a large number of results using non-uniform sampling to speed up matrix computations have appeared (Achlioptas and McSherry, 2007; Deshpande *et al.*, 2006; Deshpande and Vempala, 2006; Drineas *et al.*, 2006a,b,c,d,e; Rudelson and Vershynin, 2007; Magen and Zouzias, 2010), some of which give relative error guarantees (Deshpande *et al.*, 2006; Deshpande and Vempala, 2006; Drineas *et al.*, 2006d,e; Magen and Zouzias, 2010).

Even more recently, Sarlos (2006) showed how random projections or “sketches” can be used to perform all these tasks efficiently, obtaining the first $o(md^2)$ algorithms when preserving the identity of the rows themselves are not important. In fact, we will find many of these techniques, together with those in Ailon and Chazelle (2006) essential to our algorithm for generating row samples ultimately leading to $o(md^2)$ algorithms based on row-sampling. From now on, we focus on row-sampling algorithms.

We start with the basic result of matrix multiplication. All other results more or less follow from here. In an independent recent work which is developed along the lines of using isoperimetric inequalities (Rudelson and Vershynin, 2007) to obtain matrix Chernoff bounds, Magen and Zouzias

(2010) show that by sampling nearly a linear number of rows, it is possible to obtain a relative error approximation to matrix multiplication. Specifically, let $A \in \mathbb{R}^{m \times d_1}$ and $B \in \mathbb{R}^{m \times d_2}$. Then, for $r = \Omega(\rho/\epsilon^2 \log(d_1 + d_2))$ (where ρ bounds the stable (or “soft”) rank of A and B – see later), there is a probability distribution over $\mathcal{I} = \{1, \dots, m\}$ such that by sampling r rows i.i.d. from \mathcal{I} , one can construct sketches \tilde{A} , \tilde{B} such that $\tilde{A}^T \tilde{B} \approx A^T B$. Specifically, with constant probability,

$$\|\tilde{A}^T \tilde{B} - A^T B\|_2 \leq \epsilon \|A\|_2 \|B\|_2.$$

The sampling distribution is relatively simple, relying only on the product of the norms of the rows in A and B . This result is applied to low rank matrix reconstruction and ℓ_2 -regression where the required sampling distribution needs knowledge of the SVD of A and B .

Our basic result for matrix multiplication is very similar to this, and we arrive at it through a different path using a non-commutative Bernstein bound. Our sampling probabilities are different. In application of our results to sparse matrix reconstruction and ℓ_2 -regression, the rows of the left singular matrix make an appearance. In Magdon-Ismail (2010), it is shown how to approximate these probabilities in $o(md^2)$ time using random projections at the expense of a poly-logarithmic factor in running times. Further refinements lead to an even more efficient algorithm Drineas *et al.* (2010). As mentioned above, we must confess that one may perform our matrix tasks more efficiently using these same random projection methods (Sarlos, 2006), however the resulting algorithms are in terms of a small number of linear combinations of all the rows. In many applications, the actual rows of A have some physical meaning and so methods based on a small number of the actual rows are of interest.

We finally mention that Magen and Zouzias (2010) also give a dimension independent bound for matrix multiplication using some stronger tools. Namely, one can get the matrix multiplication approximation in the spectral norm using $r = \Omega(\rho/\epsilon^2 \log(\rho/\epsilon^2))$. In practice, it is not clear which bound is better, since there is now an additional factor of $1/\epsilon^2$ inside the logarithm.

1.1 Basic Notation

Before we can state the results in concrete form, we need some preliminary conventions. In general, $\epsilon \in (0, 1)$ will be an error tolerance parameter; $\beta \in (0, 1]$ is a parameter used to scale probabilities; and, $c, c' > 0$ are generic constants whose value may vary even within different lines of the same derivation. Let $\mathbf{e}_1, \dots, \mathbf{e}_m$ be the standard basis vectors in \mathbb{R}^m . Let $A \in \mathbb{R}^{m \times d}$ denote an arbitrary matrix which represents m points in \mathbb{R}^d . In general, we might represent a matrix such as A (roman, uppercase) by a set of vectors $\mathbf{a}_1, \dots, \mathbf{a}_m \in \mathbb{R}^d$ (bold, lowercase), so that $A^T = [\mathbf{a}_1 \ \mathbf{a}_2 \ \dots \ \mathbf{a}_m]$; similarly, for a vector \mathbf{y} , $\mathbf{y}^T = [y_1, \dots, y_m]$. Note that \mathbf{a}_t is the t^{th} row of A , which we may also refer to by $A_{(t)}$; similarly, we may refer to the t^{th} column as $A^{(t)}$. Let $\text{rank}(A) \leq \min\{m, d\}$ be the rank of A ; typically $m \gg d$ and for concreteness, we will assume that $\text{rank}(A) = d$ (all the results easily generalize to $\text{rank}(A) < d$). For matrices, we will use the spectral norm, $\|\cdot\|$; on occasion, we will use the Frobenius norm, $\|\cdot\|_F$. For vectors, $\|\cdot\|_F = \|\cdot\|$ (the standard Euclidean norm). The stable, or “soft” rank, $\rho(A) = \|A\|_F^2 / \|A\|^2 \leq \text{rank}(A)$.

The singular value decomposition (SVD) of A is

$$A = U_A S_A V_A^T.$$

where U_A is an $m \times d$ set of columns which are an orthonormal basis for the column space in A ; S_A is a $d \times d$ positive diagonal matrix of singular values, and V is a $d \times d$ orthogonal matrix. We refer to the singular values of A (the diagonal entries in S_A) by $\sigma_i(A)$. We will call a matrix with orthonormal columns an orthonormal matrix; an orthogonal matrix is a square orthonormal

matrix. In particular, $U_A^T U_A = V_A^T V_A = V_A V_A^T = I_{d \times d}$. It is possible to extend U_A to a full orthonormal basis of \mathbb{R}^m , $[U_A, U_A^\perp]$.

The SVD is important for a number of reasons. The projection of the columns of A onto the k left singular vectors with top k singular values gives the best rank- k approximation to A in the spectral and Frobenius norms. The solution to the linear regression problem is also intimately related to the SVD. In particular, consider the following minimization problem which is minimized at \mathbf{w}^* :

$$Z^* = \min_{\mathbf{w}} \|\mathbf{A}\mathbf{w} - \mathbf{y}\|^2.$$

It is known (Golub and van Loan, 1996) that $Z^* = \|\mathbf{U}_A^\perp (\mathbf{U}_A^\perp)^T \mathbf{y}\|^2$, and $\mathbf{w}^* = \mathbf{V}_A \mathbf{S}_A^{-1} \mathbf{U}_A^T \mathbf{y}$.

Row-Sampling Matrices Our focus is algorithms based on row-sampling. A *row-sampling matrix* $\mathbf{Q} \in \mathbb{R}^{r \times m}$ samples r rows of A to form $\tilde{A} = \mathbf{Q}A$:

$$\mathbf{Q} = \begin{bmatrix} \mathbf{r}_1^T \\ \vdots \\ \mathbf{r}_r^T \end{bmatrix}, \quad \tilde{A} = \mathbf{Q}A = \begin{bmatrix} \mathbf{r}_1^T A \\ \vdots \\ \mathbf{r}_r^T A \end{bmatrix} = \begin{bmatrix} \lambda_{t_1} \mathbf{a}_{t_1}^T \\ \vdots \\ \lambda_{t_r} \mathbf{a}_{t_r}^T \end{bmatrix},$$

where $\mathbf{r}_j = \lambda_{t_j} \mathbf{e}_{t_j}$; it is easy to verify that the row $\mathbf{r}_j^T A$ samples the t_j^{th} row of A and rescales it. We are interested in random sampling matrices where each \mathbf{r}_j is i.i.d. according to some distribution. Define a set of sampling probabilities p_1, \dots, p_m , with $p_i \geq 0$ and $\sum_{i=1}^m p_i = 1$; then $\mathbf{r}_j = \mathbf{e}_t / \sqrt{r p_t}$ with probability p_t . Note that the scaling is also related to the sampling probabilities in all the algorithms we consider. We can write $\mathbf{Q}^T \mathbf{Q}$ as the sum of r independently sampled matrices,

$$\mathbf{Q}^T \mathbf{Q} = \frac{1}{r} \sum_{j=1}^r \mathbf{r}_j \mathbf{r}_j^T$$

where $\mathbf{r}_j \mathbf{r}_j^T$ is a diagonal matrix with only one non-zero diagonal entry; the t^{th} diagonal entry is equal to $1/p_t$ with probability p_t . Thus, by construction, for any set of non-zero sampling probabilities, $\mathbb{E}[\mathbf{r}_j \mathbf{r}_j^T] = \mathbf{I}_{m \times m}$. Since we are averaging r independent copies, it is reasonable to expect a concentration around the mean, with respect to r , and so in some sense, $\mathbf{Q}^T \mathbf{Q}$ essentially behaves like the identity.

1.2 Statement of Results

The two main results relate to how orthonormal subspaces behave with respect to the row-sampling. These are discussed more thoroughly in Section 3, but we state them here summarily.

Theorem 1 (Symmetric Orthonormal Subspace Sampling). *Let $\mathbf{U} \in \mathbb{R}^{m \times d}$ be orthonormal, and $\mathbf{S} \in \mathbb{R}^{d \times d}$ be positive diagonal. Assume the row-sampling probabilities p_t satisfy*

$$p_t \geq \beta \frac{\mathbf{u}_t^T \mathbf{S}^2 \mathbf{u}_t}{\text{trace}(\mathbf{S}^2)}.$$

Then, if $r \geq (4\rho(\mathbf{S})/\beta\epsilon^2) \ln \frac{2d}{\delta}$, with probability at least $1 - \delta$,

$$\|\mathbf{S}^2 - \mathbf{S} \mathbf{U}^T \mathbf{Q}^T \mathbf{Q} \mathbf{U} \mathbf{S}\| \leq \epsilon \|\mathbf{S}\|^2$$

We also have an asymmetric version of Theorem 1, which is actually obtained through an application of Theorem 1 to a composite matrix.

Theorem 2 (Asymmetric Orthonormal Subspace Sampling). *Let $W \in \mathbb{R}^{m \times d_1}$, $V \in \mathbb{R}^{m \times d_2}$ be orthonormal, and let $S_1 \in \mathbb{R}^{d_1 \times d_1}$ and $S_2 \in \mathbb{R}^{d_2 \times d_2}$ be two positive diagonal matrices; let $\rho_i = \rho(S_i)$. Consider row sampling probabilities*

$$p_t \geq \beta \frac{\frac{1}{\|S_1\|^2} \mathbf{w}_t^T S_1^2 \mathbf{w}_t + \frac{1}{\|S_2\|^2} \mathbf{v}_t^T S_2^2 \mathbf{v}_t}{\rho_1 + \rho_2}.$$

If $r \geq (8(\rho_1 + \rho_2)/\beta\epsilon^2) \ln \frac{2(d_1+d_2)}{\delta}$, then with probability at least $1 - \delta$,

$$\|S_1 W^T V S_2 - S_1 W^T Q^T Q V S_2\| \leq \epsilon \|S_1\| \|S_2\|$$

We note that these row sampling probabilities are not the usual product row sampling probabilities one uses for matrix multiplication as in Drineas *et al.* (2006a). Computing the probabilities requires knowledge of the spectral norms of S_i . Here, S_i are given diagonal matrices, so it is easy to compute $\|S_i\|$. In the application of these results to matrix multiplication, the spectral norm of the input matrices will appear. We will show how to handle this issue later. As a byproduct, we will give an efficient algorithm to obtain a relative error approximation to $\|A\|$ based on row sampling and the power-iteration, which improves upon Woolfe *et al.* (2008); Kuczyński and Woźniakowski (1989).

We now give some applications of these orthonormal subspace sampling results.

Theorem 3 (Matrix Multiplication in Spectral Norm). *Let $A \in \mathbb{R}^{m \times d_1}$ and $B \in \mathbb{R}^{m \times d_2}$ have rescaled rows $\hat{\mathbf{a}}_t = \mathbf{a}_t / \|A\|$ and $\hat{\mathbf{b}}_t = \mathbf{b}_t / \|B\|$ respectively. Let ρ_A (resp. ρ_B) be the stable rank of A (resp. B). Obtain a sampling matrix $Q \in \mathbb{R}^{r \times m}$ using row-sampling probabilities p_t satisfying*

$$p_t \geq \beta \frac{\hat{\mathbf{a}}_t^T \hat{\mathbf{a}}_t + \hat{\mathbf{b}}_t^T \hat{\mathbf{b}}_t}{\sum_{t=1}^m \hat{\mathbf{a}}_t^T \hat{\mathbf{a}}_t + \hat{\mathbf{b}}_t^T \hat{\mathbf{b}}_t} = \beta \frac{\hat{\mathbf{a}}_t^T \hat{\mathbf{a}}_t + \hat{\mathbf{b}}_t^T \hat{\mathbf{b}}_t}{\rho_A + \rho_B}.$$

Then, if $r \geq \frac{8(\rho_A + \rho_B)}{\beta\epsilon^2} \ln \frac{2(d_1+d_2)}{\delta}$, with probability at least $1 - \delta$,

$$\|A^T B - \tilde{A}^T \tilde{B}\| \leq \epsilon \|A\| \|B\|.$$

The sampling probabilities depend on $\|A\|^2$ and $\|B\|^2$. It is possible to get a constant factor approximation to $\|A\|^2$ (and similarly $\|B\|^2$) with high probability. We summarize the idea here, the details are given in Section 7, Theorem 25. First sample $\tilde{A} = QA$ according to probabilities $p_t = \mathbf{a}_t^2 / \|A\|_F^2$. These probabilities are easy to compute in $O(md_1)$. By an application of the symmetric subspace sampling theorem (see Theorem 17), if $r \geq (4\rho_A/\epsilon^2) \ln \frac{2d_1}{\delta}$, then with probability at least $1 - \delta$,

$$(1 - \epsilon) \|A\|^2 \leq \|\tilde{A}^T \tilde{A}\| \leq (1 + \epsilon) \|A\|^2.$$

We now run $\Omega(\ln \frac{d_1}{\delta})$ power iterations starting from a random isotropic vector to estimate the spectral norm of $\tilde{A}^T \tilde{A}$. The efficiency is $O(md_1 + \rho_A d_1 / \epsilon^2 \ln^2(\frac{d_1}{\delta}))$.

Theorem 4 (Sparse Row-Based Matrix Reconstruction). *Let A have the SVD representation $A = USV^T$, and consider row-sampling probabilities p_t satisfying $p_t \geq \frac{\beta}{d} \mathbf{u}_t^T \mathbf{u}_t$. Then, if $r \geq (4(d - \beta)/\beta\epsilon^2) \ln \frac{2d}{\delta}$, with probability at least $1 - \delta$,*

$$\|A - A\tilde{\Pi}_k\| \leq \left(\frac{1 + \epsilon}{1 - \epsilon} \right)^{1/2} \|A - A_k\|,$$

for $k = 1, \dots, d$, where $\tilde{\Pi}_k$ projects onto the top k right singular vectors of \tilde{A} .

It is possible to obtain relative approximations to the sampling probabilities according to the rows of the left singular matrix (the leverage scores), but that goes beyond the scope of this work Magdon-Ismail (2010); Drineas *et al.* (2010)

Theorem 5 (Relative Error ℓ_2 Regression). *Let $A \in \mathbb{R}^{m \times d}$ have the SVD representation $A = USV^T$, and let $\mathbf{y} \in \mathbb{R}^m$. Let $\mathbf{x}^* = A^+ \mathbf{y}$ be the optimal regression with residual $\boldsymbol{\epsilon} = \mathbf{y} - A\mathbf{x}^* = \mathbf{y} - AA^+ \mathbf{y}$. Assume the sampling probabilities p_t satisfy*

$$p_t \geq \beta \left(\frac{\mathbf{u}_t^2}{d} + \frac{(\mathbf{u}_t^2 + \frac{\epsilon_t^2}{\boldsymbol{\epsilon}^T \boldsymbol{\epsilon}})}{d+1} + \frac{\epsilon_t^2}{\boldsymbol{\epsilon}^T \boldsymbol{\epsilon}} \right)$$

For $r \geq (8(d+1)/\beta\epsilon^2) \ln \frac{2(d+1)}{\delta}$, let $\hat{\mathbf{x}} = (QA)^+ Q\mathbf{y}$ be the approximate regression. Then, with probability at least $1 - 3\delta$,

$$\|A\hat{\mathbf{x}} - \mathbf{y}\| \leq \left(1 + \epsilon + \epsilon \sqrt{\frac{1+\epsilon}{1-\epsilon}} \right) \|A\mathbf{x}^* - \mathbf{y}\|.$$

In addition to sampling according to \mathbf{u}_t^2 we also need the residual vector $\boldsymbol{\epsilon} = \mathbf{y} - AA^+ \mathbf{y}$. Unfortunately, we have not yet found an efficient way to get a good approximation (in some form of relative error) to this residual vector.

1.3 Paper Outline

Next we describe some probabilistic tail inequalities which will be useful. We continue with the sampling lemmas for orthonormal matrices, followed by the applications to matrix multiplication, matrix reconstruction and ℓ_2 -regression. Finally, we discuss the algorithm for approximating the spectral norm based on sampling and the power iteration.

2 Probabilistic Tail Inequalities

Since all our arguments involve high probability results, our main bounding tools will be probability tail inequalities. First, let X_1, \dots, X_n be independent random variables with $\mathbb{E}[X_i] = 0$ and $|\mathbf{X}_i| \leq \gamma$; let $Z_n = \frac{1}{n} \sum_{i=1}^n X_i$. Chernoff, and later Hoeffding gave the bound

Theorem 6 (Chernoff (1952); Hoeffding (1963)). $\mathbb{P}[|Z_n| > \epsilon] \leq 2e^{-n\epsilon^2/2\gamma^2}$.

If in addition one can bound the variance, $\mathbb{E}[X_i^2] \leq s^2$, then we have Bernstein's bound:

Theorem 7 (Bernstein (1924)). $\mathbb{P}[|Z_n| \geq \epsilon] \leq 2e^{-n\epsilon^2/(2s^2+2\gamma\epsilon/3)}$.

Note that when $\epsilon \leq 3s^2/\gamma$, we can simplify the Bernstein bound to $\mathbb{P}[|Z_n| \geq \epsilon] \leq 2e^{-n\epsilon^2/4s^2}$, which is considerably simpler and only involves the variance. The non-commutative versions of these bounds, which extend these inequalities to matrix valued random variables can also be deduced. Let X_1, \dots, X_n be independent copies of a symmetric random matrix X , with $\mathbb{E}[X] = \mathbf{0}$, and suppose that $\|X\|_2 \leq \gamma$; let $Z_n = \frac{1}{n} \sum_{i=1}^n X_i$. Ahlswede and Winter (2002) gave the fundamental extension of the exponentiation trick for computing Chernoff bounds of scalar random variables to matrix valued random variables (for a simplified proof, see Wigderson and Xiao (2008)):

$$\mathbb{P}[\|Z_n\|_2 > \epsilon] \leq \inf_t 2de^{-n\epsilon t/\gamma} \mathbb{E} \|e^{tX/\gamma}\|_2^n. \quad (1)$$

By standard optimization of this bound, one readily obtains the non-commutative tail inequality

Theorem 8 (Ahlsvede and Winter (2002)). $\mathbb{P}[\|Z_n\|_2 > \epsilon] \leq 2de^{-n\epsilon^2/4\gamma^2}$.

Proof. The statement is trivial if $\epsilon \geq \gamma$, so assume $\epsilon < \gamma$. The lemma follows from (1) and the following sequence after setting $t = \epsilon/2\gamma \leq \frac{1}{2}$:

$$\|\mathbb{E}[e^{tX/\gamma}]\|_2 \stackrel{(a)}{\leq} 1 + \sum_{\ell=2}^{\infty} \frac{t^\ell}{\ell!} \mathbb{E}[\|(X/\gamma)^\ell\|_2] \stackrel{(b)}{\leq} 1 + t^2 \leq e^{t^2}, \quad (2)$$

where (a) follows from $\mathbb{E}[X] = 0$, the triangle inequality and $\|\mathbb{E}[\cdot]\|_2 \leq \mathbb{E}[\|\cdot\|_2]$; (b) follows because $\|X/\gamma\|_2 \leq 1$ and $t \leq \frac{1}{2}$. ■

(We have stated a simplified version of the bound, without taking care to optimize the constants.) As mentioned in Gross *et al.* (2009), one can obtain a non-commuting version of Bernstein's inequality in a similar fashion using (1). Assume that $\|\mathbb{E} X^T X\|_2 \leq s^2$. By adapting the standard Bernstein bounding argument to matrices, we have

Lemma 9. *For symmetric X , $\|\mathbb{E}[e^{tX/\gamma}]\|_2 \leq \exp\left(\frac{s^2}{\gamma^2}(e^t - 1 - t)\right)$.*

Proof. As in (2), but using submultiplicativity, we first bound $\|\mathbb{E}[X^\ell]\|_2 \leq s^2\gamma^{\ell-2}$:

$$\begin{aligned} \|\mathbb{E}[X^\ell]\|_2 &= \max_{\|\mathbf{u}\|=1} \left\| \int dX p(X) X^\ell \mathbf{u} \right\| \\ &= \max_{\|\mathbf{u}\|=1} \left\| \int dX p(X) \frac{\|X^{\ell-2}\mathbf{u}\| X^2 X^{\ell-2}\mathbf{u}}{\|X^{\ell-2}\mathbf{u}\|} \right\| \\ &\leq \gamma^{\ell-2} \max_{\|\mathbf{w}\|=1} \left\| \int dX p(X) X^2 \mathbf{w} \right\| \\ &= \gamma^{\ell-2} \|\mathbb{E}[X^2]\|_2 \leq s^2\gamma^{\ell-2}. \end{aligned}$$

To conclude, we use the triangle inequality to bound as follows:

$$\|\mathbb{E}[e^{tX/\gamma}]\|_2 = \left\| \mathbf{I} + \sum_{\ell=2}^{\infty} \frac{t^\ell}{\gamma^\ell \ell!} \mathbb{E}[X^\ell] \right\|_2 \leq 1 + \frac{s^2}{\gamma^2} \sum_{\ell=2}^{\infty} \frac{t^\ell}{\ell!} = 1 + \frac{s^2}{\gamma^2}(e^t - 1 - t) \leq \exp\left(\frac{s^2}{\gamma^2}(e^t - 1 - t)\right). \quad \blacksquare$$

Using Lemma 9 in (1) with $t = \ln(1 + \epsilon\gamma/s^2)$, and using $(1+x)\ln(1 + \frac{1}{x}) - 1 \geq \frac{1}{2x+2/3}$, we obtain the following result.

Theorem 10 (Non-commutative Bernstein). $\mathbb{P}[\|Z_n\|_2 > \epsilon] \leq 2de^{-n\epsilon^2/(2s^2+2\gamma\epsilon/3)}$.

Gross *et al.* (2009) gives a simpler version of this non-commutative Bernstein inequality. If $X \in \mathbb{R}^{d_1 \times d_2}$ is not symmetric, then by considering

$$\begin{bmatrix} \mathbf{0}_{d_1 \times d_1} & X \\ X^T & \mathbf{0}_{d_2 \times d_2} \end{bmatrix},$$

one can get a non-symmetric version of the non-commutative Chernoff and Bernstein bounds,

Theorem 11 (Recht (2009)). $\mathbb{P}[\|Z_n\|_2 > \epsilon] \leq (d_1 + d_2)e^{-n\epsilon^2/(2s^2+2\gamma\epsilon/3)}$.

For most of our purposes, we will only need the symmetric version; again, if $\epsilon \leq 3s^2/\gamma$, then we have the much simpler bound $\mathbb{P}[\|Z_n\|_2 > \epsilon] \leq 2de^{-n\epsilon^2/4s^2}$.

3 Orthonormal Sampling Lemmas

Let $U \in \mathbb{R}^{m \times d}$ be an orthonormal matrix, and let $S \in \mathbb{R}^{d \times d}$ be a diagonal matrix. We are interested in the product $US \in \mathbb{R}^{m \times d}$; US is the matrix with columns $U^{(i)}S_{ii}$. Without loss of generality, we can assume that S is positive by flipping the signs of the appropriate columns of U . The row-representation of U is $U^T = [\mathbf{u}_1, \dots, \mathbf{u}_m]$; we consider the row sampling probabilities

$$p_t \geq \beta \frac{\mathbf{u}_t^T S^2 \mathbf{u}_t}{\text{trace}(S^2)}. \quad (3)$$

Since $U^T U = I_{d \times d}$, one can verify that $\text{trace}(S^2) = \sum_t \mathbf{u}_t^T S^2 \mathbf{u}_t$ is the correct normalization.

Lemma 12 (Symmetric Subspace Sampling Lemma).

$$\begin{aligned} \mathbb{P}[\|S^2 - SU^T Q^T Q US\| > \epsilon \|S\|^2] &\leq 2d \cdot \exp\left(\frac{-r\epsilon^2}{2(\rho/\beta - \kappa^{-4} + \epsilon(\rho/\beta - \kappa^{-2})/3)}\right), \\ &\leq 2d \cdot \exp\left(\frac{-r\beta\epsilon^2}{4\rho}\right), \end{aligned}$$

where ρ is the numerical (stable) rank of S , $\rho(S) = \|S\|_F^2 / \|S\|^2$, and $\kappa(S) = \sigma_{\max}(S) / \sigma_{\min}(S)$ is the condition number.

Remarks. The stable rank $\rho \leq d$ measures the effective dimension of the matrix. The condition number $\kappa \geq 1$, hence the simpler version of the bound, which is valid for $\epsilon \leq 3$. It immediately follows that if $r \geq (4\rho/\beta\epsilon^2) \ln \frac{2d}{\delta}$, then with probability at least $1 - \delta$,

$$\|S^2 - SU^T Q^T Q US\| \leq \epsilon \|S\|^2$$

An important special case is when $S = I_{d \times d}$, in which case $\rho = d$, $\kappa = 1$ and $\|S\| = 1$.

Corollary 13. For sampling probabilities $p_t \geq \frac{\beta}{d} \mathbf{u}_t^T \mathbf{u}_t$,

$$\mathbb{P}[\|I - U^T Q^T Q U\| > \epsilon] \leq 2d \cdot \exp\left(\frac{-\beta r \epsilon^2}{4(d - \beta)}\right),$$

Proof. (of Lemma 12) Note that $U^T Q^T Q U = \frac{1}{r} \sum_{i=1}^r \mathbf{u}_{t_i} \mathbf{u}_{t_i}^T / p_{t_i}$, where $t_i \in [1, m]$ is chosen according to the probability p_{t_i} . It follows that

$$S^2 - SU^T Q^T Q US = \frac{1}{r} \sum_{i=1}^r S^2 - \frac{1}{p_{t_i}} S \mathbf{u}_{t_i} \mathbf{u}_{t_i}^T S = \frac{1}{r} \sum_{i=1}^r X_i,$$

where X_i are independent copies of a matrix-random variable $X \sim S^2 - S \mathbf{u} \mathbf{u}^T S / p$. We prove the following three claims:

- (i) $\mathbb{E}[X] = \mathbf{0}$;
- (ii) $\|X\| \leq \|S\|^2(\rho/\beta - \kappa^{-2})$;
- (ii) $\|\mathbb{E} X^T X\| \leq \|S\|^4(\rho/\beta - \kappa^{-4})$.

The Lemma follows from the non-commutative Bernstein bound with ϵ replaced by $\epsilon \|S\|^2$. To prove (i), note that $\mathbb{E}[X] = S^2 - S \mathbb{E}[\mathbf{u} \mathbf{u}^T / p] S = S^2 - S (\sum_{t=1}^m \mathbf{u}_t \mathbf{u}_t^T) S = \mathbf{0}$, because $\sum_{t=1}^m \mathbf{u}_t \mathbf{u}_t^T = U^T U = I_{d \times d}$.

To prove (ii), let \mathbf{z} be an arbitrary unit vector and consider

$$\mathbf{z}^T \mathbf{X} \mathbf{z} = \mathbf{z}^T \mathbf{S}^2 \mathbf{z} - \frac{1}{p} (\mathbf{z}^T \mathbf{S} \mathbf{u})^2.$$

It follows that $\mathbf{z}^T \mathbf{X} \mathbf{z} \leq \|\mathbf{S}\|^2$. To get a lower bound, we use $p \geq \beta \mathbf{u}^T \mathbf{S}^2 \mathbf{u} / \text{trace}(\mathbf{S}^2)$:

$$\begin{aligned} \mathbf{z}^T \mathbf{X} \mathbf{z} &\geq \mathbf{z}^T \mathbf{S}^2 \mathbf{z} - \frac{\text{trace}(\mathbf{S}^2)}{\beta} \frac{(\mathbf{z}^T \mathbf{S} \mathbf{u})^2}{\mathbf{u}^T \mathbf{S}^2 \mathbf{u}}, \\ &\stackrel{(a)}{\geq} \|\mathbf{S}\|^2 \left(\frac{\sigma_{\min}^2(\mathbf{S})}{\|\mathbf{S}\|^2} - \frac{\text{trace}(\mathbf{S}^2)}{\beta \|\mathbf{S}\|^2} \right), \\ &= \|\mathbf{S}\|^2 \left(\frac{1}{\kappa^2} - \frac{\rho}{\beta} \right). \end{aligned}$$

(a) follows because: by definition of σ_{\min} , the minimum of the first term is σ_{\min}^2 ; and, by Cauchy-Schwarz, $(\mathbf{z}^T \mathbf{S} \mathbf{u})^2 \leq (\mathbf{z}^T \mathbf{z})(\mathbf{u}^T \mathbf{S}^2 \mathbf{u})$. Since $\beta \leq 1$, $\rho/\beta - \kappa^{-2} \geq 1$ (for $d > 1$), and so $|\mathbf{z}^T \mathbf{X} \mathbf{z}| \leq \|\mathbf{S}\|^2 (\rho/\beta - \kappa^{-2})$, from which (ii) follows.

To prove (iii), first note that

$$\begin{aligned} \mathbb{E}[\mathbf{X}^T \mathbf{X}] &= \mathbf{S}^4 - \mathbf{S}^3 \mathbb{E}[\mathbf{u} \mathbf{u}^T / p] \mathbf{S} - \mathbf{S} \mathbb{E}[\mathbf{u} \mathbf{u}^T / p] \mathbf{S}^3 + \mathbf{S} \mathbb{E}[\mathbf{u} \mathbf{u}^T \mathbf{S}^2 \mathbf{u} \mathbf{u}^T / p^2] \mathbf{S}, \\ &\stackrel{(a)}{=} \mathbf{S} \left(\sum_{t=1}^m \frac{1}{p_t} \mathbf{u}_t \mathbf{u}_t^T \mathbf{S}^2 \mathbf{u}_t \mathbf{u}_t^T \right) \mathbf{S} - \mathbf{S}^4. \end{aligned}$$

(a) follows because $\mathbb{E}[\mathbf{u} \mathbf{u}^T / p] = \mathbf{I}$. Thus, for an arbitrary unit \mathbf{z} , we have

$$\begin{aligned} \mathbf{z}^T \mathbb{E}[\mathbf{X}^T \mathbf{X}] \mathbf{z} &= \sum_{t=1}^m \frac{1}{p_t} (\mathbf{z}^T \mathbf{S} \mathbf{u}_t \mathbf{u}_t^T \mathbf{S} \mathbf{z}) \mathbf{u}_t^T \mathbf{S}^2 \mathbf{u}_t - \mathbf{z}^T \mathbf{S}^4 \mathbf{z}, \\ &\stackrel{(a)}{\leq} \frac{\text{trace}(\mathbf{S}^2)}{\beta} \mathbf{z}^T \mathbf{S} \left(\sum_{t=1}^m \mathbf{u}_t \mathbf{u}_t^T \right) \mathbf{S} \mathbf{z} - \mathbf{z}^T \mathbf{S}^4 \mathbf{z}, \\ &\stackrel{(b)}{=} \|\mathbf{S}\|^4 \left(\frac{\text{trace}(\mathbf{S}^2)}{\beta \|\mathbf{S}\|^2} \frac{\mathbf{z}^T \mathbf{S}^2 \mathbf{z}}{\|\mathbf{S}\|^2} - \frac{\mathbf{z}^T \mathbf{S}^4 \mathbf{z}}{\|\mathbf{S}\|^4} \right), \\ &\leq \|\mathbf{S}\|^4 \left(\frac{\text{trace}(\mathbf{S}^2)}{\beta \|\mathbf{S}\|^2} - \frac{\sigma_{\min}^4}{\|\mathbf{S}\|^4} \right). \end{aligned}$$

(a) follows from $p_t \geq \beta \mathbf{u}_t^T \mathbf{S}^2 \mathbf{u}_t / \text{trace}(\mathbf{S}^2)$; (b) follows from $\mathbf{U}^T \mathbf{U} = \sum_{t=1}^m \mathbf{u}_t \mathbf{u}_t^T = \mathbf{I}_{d \times d}$. Thus, $|\mathbf{z}^T \mathbb{E}[\mathbf{X}^T \mathbf{X}] \mathbf{z}| \leq \|\mathbf{S}\|^4 (\rho/\beta - \kappa^{-4})$, from which (iii) follows. \blacksquare

For the general case, consider two orthonormal matrices $\mathbf{W} \in \mathbb{R}^{m \times d_1}$, $\mathbf{V} \in \mathbb{R}^{m \times d_2}$, and two positive diagonal matrices $\mathbf{S}_1 \in \mathbb{R}^{d_1 \times d_1}$ and $\mathbf{S}_2 \in \mathbb{R}^{d_2 \times d_2}$. We consider the product $\mathbf{S}_1 \mathbf{W}^T \mathbf{V} \mathbf{S}_2$, which is approximated by the sampled product $\mathbf{S}_1 \mathbf{W}^T \mathbf{Q}^T \mathbf{Q} \mathbf{V} \mathbf{S}_2$. Consider the sampling probabilities

$$p_t \geq \beta \frac{(\mathbf{u}_t^T \mathbf{S}_1^2 \mathbf{u}_t)^{1/2} (\mathbf{v}_t^T \mathbf{S}_2^2 \mathbf{v}_t)^{1/2}}{\sum_{t=1}^m (\mathbf{u}_t^T \mathbf{S}_1^2 \mathbf{u}_t)^{1/2} (\mathbf{v}_t^T \mathbf{S}_2^2 \mathbf{v}_t)^{1/2}} \geq \beta \frac{(\mathbf{u}_t^T \mathbf{S}_1^2 \mathbf{u}_t)^{1/2} (\mathbf{v}_t^T \mathbf{S}_2^2 \mathbf{v}_t)^{1/2}}{\sqrt{\text{trace}(\mathbf{S}_1^2) \text{trace}(\mathbf{S}_2^2)}},$$

where the last inequality follows from Cauchy-Schwarz. Since $\|\mathbf{A}\|_F = \sqrt{\rho(\mathbf{A})} \|\mathbf{A}\| \geq \|\mathbf{A}\|$, any bound for the Frobenius norm can be converted into a bound for the spectral norm. Using the

Frobenius norm bounds in Drineas *et al.* (2006a) (using a simplified form for the bound), one immediately has:

$$\mathbb{P} [\|S_1 W^T V S_2 - S_1 W^T Q^T Q V S_2\| > \epsilon \|S_1\| \|S_2\|] \leq \exp \left(\frac{-r \beta^2 \epsilon^2}{16 \rho_1 \rho_2} \right), \quad (4)$$

where $\rho_1 = \rho(S_1)$ and $\rho_2 = \rho(S_2)$. Alternatively, if $r \geq (16 \rho_1 \rho_2 / \beta^2 \epsilon^2) \ln \frac{1}{\delta}$, then

$$\|S_1 W^T V S_2 - S_1 W^T Q^T Q V S_2\| \leq \epsilon \|S_1\| \|S_2\|.$$

The dependence on the stable ranks and β is quadratic. Applying this bound to the situation in Lemma 12 would give an inferior bound. The intuition behind the improvement is that the sampling is isotropic, and so will not favor any particular direction. One can therefore guess that all the singular values are approximately equal and so the Frobenius norm bound on the spectral norm will be loose by a factor of $\sqrt{\rho}$; and, indeed this is what comes out in the closer analysis. As a application of Lemma 12, we can get a better result for the asymmetric case.

Lemma 14. *Let $W \in \mathbb{R}^{m \times d_1}$, $V \in \mathbb{R}^{m \times d_2}$ be orthonormal, and let $S_1 \in \mathbb{R}^{d_1 \times d_1}$ and $S_2 \in \mathbb{R}^{d_2 \times d_2}$ be two positive diagonal matrices. Consider row sampling probabilities*

$$p_t \geq \beta \frac{\frac{1}{\|S_1\|^2} \mathbf{w}_t^T S_1^2 \mathbf{w}_t + \frac{1}{\|S_2\|^2} \mathbf{v}_t^T S_2^2 \mathbf{v}_t}{\rho_1 + \rho_2}.$$

If $r \geq (8(\rho_1 + \rho_2) / \beta \epsilon^2) \ln \frac{2(d_1 + d_2)}{\delta}$, then with probability at least $1 - \delta$,

$$\|S_1 W^T V S_2 - S_1 W^T Q^T Q V S_2\| \leq \epsilon \|S_1\| \|S_2\|$$

For the special case that $S_1 = I_{d_1 \times d_1}$ and $S_2 = I_{d_2 \times d_2}$, the sampling probabilities simplify to

$$p_t \geq \beta \frac{\mathbf{w}_t^T \mathbf{w}_t + \mathbf{v}_t^T \mathbf{v}_t}{d_1 + d_2},$$

Corollary 15. *If $r \geq (8(d_1 + d_2) / \beta \epsilon^2) \ln \frac{2(d_1 + d_2)}{\delta}$, then with probability at least $1 - \delta$,*

$$\|W^T V - W^T Q^T Q V\| \leq \epsilon.$$

Proof. (of Lemma 14) By homogeneity, we can without loss of generality assume that $\|S_1\| = \|S_2\| = 1$, and let¹ $Z = [W S_1 \ V S_2]$. An elementary lemma which we will find useful is

Lemma 16. *For any matrix $A = [A_1 \ A_2]$,*

$$\max\{\|A_1\|, \|A_2\|\} \leq \|A\| \leq \sqrt{\|A_1\|^2 + \|A_2\|^2}.$$

The left inequality is saturated when A_1 and A_2 are orthogonal ($A_1^T A_2 = \mathbf{0}$), and the right inequality is saturated when $A_1 = A_2$. By repeatedly applying Lemma 16 one can see that $\|A\|$ is at least the spectral norm of any submatrix. Introduce the SVD of Z ,

$$Z = [W S_1 \ V S_2] = U S V_Z^T.$$

¹The general case would have been $Z = \left[\frac{1}{\|S_1\|} W S_1 \ \frac{1}{\|S_2\|} V S_2 \right]$.

We now use the row sampling probabilities according to US from (3),

$$p_t \geq \beta \frac{\mathbf{u}_t^T \mathbf{S}^2 \mathbf{u}_t}{\text{trace}(\mathbf{S}^2)}.$$

We may interpret the sampling probabilities as follows. Let \mathbf{z}_t be a row of \mathbf{Z} , the concatenation of two rows in $\mathbf{W}\mathbf{S}_1$ and $\mathbf{V}\mathbf{S}_2$: $\mathbf{z}_t^T = [\mathbf{w}_t^T \mathbf{S}_1 \quad \mathbf{v}_t^T \mathbf{S}_2]$. We also have that $\mathbf{z}_t^T = \mathbf{u}_t^T \mathbf{S} \mathbf{V}_Z^T$. Hence,

$$\mathbf{u}_t^T \mathbf{S}^2 \mathbf{u}_t = \mathbf{u}_t^T \mathbf{S} \mathbf{V}_Z^T \mathbf{V}_Z \mathbf{S} \mathbf{u}_t = \mathbf{z}_t^T \mathbf{z}_t = \mathbf{w}_t^T \mathbf{S}_1^2 \mathbf{w}_t + \mathbf{v}_t^T \mathbf{S}_2^2 \mathbf{v}_t.$$

These are exactly the probabilities as claimed in the statement of the lemma (modulo the rescaling).

Applying Lemma 12: if $r \geq (4\rho/\beta\epsilon^2) \ln \frac{2\text{rank}(\mathbf{U})}{\delta}$, then with probability at least $1 - \delta$,

$$\|\mathbf{S}^2 - \mathbf{S} \mathbf{U}^T \mathbf{Q}^T \mathbf{Q} \mathbf{U} \mathbf{S}\| \leq \epsilon \|\mathbf{S}\|^2 \leq \epsilon \sqrt{\|\mathbf{S}_1\|^2 + \|\mathbf{S}_2\|^2} = \epsilon \sqrt{2},$$

where the second inequality follows from Lemma 16. Since $\mathbf{Z}\mathbf{V} = \mathbf{U}\mathbf{S}$,

$$\|\mathbf{Z}^T \mathbf{Z} - \mathbf{Z}^T \mathbf{Q}^T \mathbf{Q} \mathbf{Z}\| = \|\mathbf{S}^2 - \mathbf{S} \mathbf{U}^T \mathbf{Q}^T \mathbf{Q} \mathbf{U} \mathbf{S}\|.$$

Further, by the construction of \mathbf{Z} ,

$$\mathbf{Z}^T \mathbf{Z} - \mathbf{Z}^T \mathbf{Q}^T \mathbf{Q} \mathbf{Z} = \begin{bmatrix} \mathbf{S}_1^2 - \mathbf{S}_1 \mathbf{W}^T \mathbf{Q}^T \mathbf{Q} \mathbf{W} \mathbf{S}_1 & \mathbf{S}_1 \mathbf{W}^T \mathbf{V} \mathbf{S}_2 - \mathbf{S}_1 \mathbf{W}^T \mathbf{Q}^T \mathbf{Q} \mathbf{V} \mathbf{S}_2 \\ \mathbf{S}_2 \mathbf{V}^T \mathbf{W} \mathbf{S}_1 - \mathbf{S}_2 \mathbf{V}^T \mathbf{Q}^T \mathbf{Q} \mathbf{W} \mathbf{S}_1 & \mathbf{S}_2^2 - \mathbf{S}_2 \mathbf{V}^T \mathbf{Q}^T \mathbf{Q} \mathbf{V} \mathbf{S}_2 \end{bmatrix}.$$

By Lemma 16, $\|\mathbf{S}_1 \mathbf{W}^T \mathbf{V} \mathbf{S}_2 - \mathbf{S}_1 \mathbf{W}^T \mathbf{Q}^T \mathbf{Q} \mathbf{V} \mathbf{S}_2\| \leq \|\mathbf{Z}^T \mathbf{Z} - \mathbf{Z}^T \mathbf{Q}^T \mathbf{Q} \mathbf{Z}\|$, and so:

$$\|\mathbf{S}_1 \mathbf{W}^T \mathbf{V} \mathbf{S}_2 - \mathbf{S}_1 \mathbf{W}^T \mathbf{Q}^T \mathbf{Q} \mathbf{V} \mathbf{S}_2\| \leq \epsilon \sqrt{2}.$$

Observe that $\text{trace}(\mathbf{S}^2) = \|\mathbf{Z}\|_F^2 = \text{trace}(\mathbf{S}_1^2) + \text{trace}(\mathbf{S}_2^2)$; further, since $\|\mathbf{S}\| \geq \max\{\|\mathbf{S}_1\|, \|\mathbf{S}_2\|\}$, we have that

$$\rho(\mathbf{S}) = \frac{\text{trace}(\mathbf{S}^2)}{\|\mathbf{S}\|^2} = \frac{\text{trace}(\mathbf{S}_1^2) + \text{trace}(\mathbf{S}_2^2)}{\|\mathbf{S}\|^2} \leq \frac{\text{trace}(\mathbf{S}_1^2)}{\|\mathbf{S}_1\|^2} + \frac{\text{trace}(\mathbf{S}_2^2)}{\|\mathbf{S}_2\|^2} = \rho_1 + \rho_2.$$

Since $\text{rank}(\mathbf{U}) \leq d_1 + d_2$, it suffices that $r \geq (4(\rho_1 + \rho_2)/\beta\epsilon^2) \ln \frac{2(d_1+d_2)}{\delta}$ to obtain error $\epsilon\sqrt{2}$; after rescaling $\epsilon' = \epsilon\sqrt{2}$, we have the result. \blacksquare

4 Sampling for Matrix Multiplication

We obtain results for matrix multiplication directly from Lemmas 12 and 14. First we consider the symmetric case, then the asymmetric case. Let $\mathbf{A} \in \mathbb{R}^{m \times d_1}$ and $\mathbf{B} \in \mathbb{R}^{m \times d_2}$. We are interested in conditions on the sampling matrix $\mathbf{Q} \in \mathbb{R}^{r \times m}$ such that $\mathbf{A}^T \mathbf{A} \approx \tilde{\mathbf{A}}^T \tilde{\mathbf{A}}$ and $\mathbf{A}^T \mathbf{B} \approx \tilde{\mathbf{A}}^T \tilde{\mathbf{B}}$, where $\tilde{\mathbf{A}} = \mathbf{Q} \mathbf{A}$ and $\tilde{\mathbf{B}} = \mathbf{Q} \mathbf{B}$. Using the SVD of \mathbf{A} ,

$$\begin{aligned} \|\mathbf{A}^T \mathbf{A} - \mathbf{A}^T \mathbf{Q}^T \mathbf{Q} \mathbf{A}\| &= \|\mathbf{V}_A \mathbf{S}_A \mathbf{U}_A^T \mathbf{U}_A \mathbf{S}_A \mathbf{V}_A^T - \mathbf{V}_A \mathbf{S}_A \mathbf{U}_A^T \mathbf{Q}^T \mathbf{Q} \mathbf{U}_A \mathbf{S}_A \mathbf{V}_A^T\|, \\ &= \|\mathbf{S}_A^2 - \mathbf{S}_A \mathbf{U}_A^T \mathbf{Q}^T \mathbf{Q} \mathbf{U}_A \mathbf{S}_A\|. \end{aligned}$$

We may now directly apply Lemma 12, with respect to the appropriate sampling probabilities. One can verify that the sampling probabilities in Lemma 12 are proportional to the squared norms of the rows of \mathbf{A} . \blacksquare

Theorem 17. Let $A \in \mathbb{R}^{m \times d_1}$ have rows \mathbf{a}_t . Obtain a sampling matrix $Q \in \mathbb{R}^{r \times m}$ using row-sampling probabilities

$$p_t \geq \beta \frac{\mathbf{a}_t^T \mathbf{a}_t}{\|A\|_F^2}.$$

Then, if $r \geq \frac{4\rho_A}{\beta\epsilon^2} \ln \frac{2d_1}{\delta}$, with probability at least $1 - \delta$,

$$\|A^T A - \tilde{A}^T \tilde{A}\| \leq \epsilon \|A\|^2.$$

Similarly, using the SVDs of A and B ,

$$\begin{aligned} \|A^T B - A^T Q^T Q B\| &= \|V_A S_A U_A^T U_B S_B V_B^T - V_A S_A U_A^T Q^T Q U_B S_B V_B^T\|, \\ &= \|S_A U_A^T U_B S_B - S_A U_A^T Q^T Q U_B S_B\|. \end{aligned}$$

We may now directly apply Lemma 14, with respect to the appropriate sampling probabilities. One can verify that the sampling probabilities in Lemma 14 are proportional to the sum of the rescaled squared norms of the rows of A and B .

Theorem 18. Let $A \in \mathbb{R}^{m \times d_1}$ and $B \in \mathbb{R}^{m \times d_2}$, have rescaled rows $\hat{\mathbf{a}}_t = \mathbf{a}_t / \|A\|$ and $\hat{\mathbf{b}}_t = \mathbf{b}_t / \|B\|$ respectively. Obtain a sampling matrix $Q \in \mathbb{R}^{r \times m}$ using row-sampling probabilities

$$p_t \geq \beta \frac{\hat{\mathbf{a}}_t^T \hat{\mathbf{a}}_t + \hat{\mathbf{b}}_t^T \hat{\mathbf{b}}_t}{\sum_{t=1}^m \hat{\mathbf{a}}_t^T \hat{\mathbf{a}}_t + \hat{\mathbf{b}}_t^T \hat{\mathbf{b}}_t} = \beta \frac{\hat{\mathbf{a}}_t^T \hat{\mathbf{a}}_t + \hat{\mathbf{b}}_t^T \hat{\mathbf{b}}_t}{\rho_A + \rho_B}.$$

Then, if $r \geq \frac{8(\rho_A + \rho_B)}{\beta\epsilon^2} \ln \frac{2(d_1 + d_2)}{\delta}$, with probability at least $1 - \delta$,

$$\|A^T B - \tilde{A}^T \tilde{B}\| \leq \epsilon \|A\| \|B\|.$$

5 Sparse Row Based Matrix Representation

Given a matrix $A = USV^T \in \mathbb{R}^{m \times d}$, the top k singular vectors, corresponding to the top k singular values give the best rank k reconstruction of A . Specifically, let $A_k = U_k S_k V_k^T$, where $U_k \in \mathbb{R}^{m \times k}$, $S_k \in \mathbb{R}^{k \times k}$ and $V_k \in \mathbb{R}^{d \times k}$; U_k and V_k correspond to the top- k left and right singular vectors. Then, $\|A - A_k\| \leq \|A - X\|$ where $X \in \mathbb{R}^{m \times d}$ ranges over all rank- k matrices. As usual, let $\tilde{A} = QA$ be the sampled, rescaled rows of A , with $\tilde{A} = \tilde{U} \tilde{S} \tilde{V}^T$, and consider the top- k right singular vectors \tilde{V}_k . Let $\tilde{\Pi}_k$ be the projection onto this top- k right singular space, and consider the rank k approximation to A obtained by projecting onto this space: $\tilde{A}_k = \tilde{A} \tilde{\Pi}_k$. The following lemma is useful for showing that \tilde{A}_k is almost (up to additive error) as good an approximation to A as one can get.

Lemma 19 (Drineas *et al.* (2006b), Rudelson and Vershynin (2007)).

$$\|A - \tilde{A}_k\|^2 \leq \|A - A_k\|^2 + 2\|A^T A - \tilde{A}^T \tilde{A}\| \leq (\|A - A_k\| + \sqrt{2}\|A^T A - \tilde{A}^T \tilde{A}\|^{1/2})^2.$$

Proof. The proof follows using standard arguments and an application of a perturbation theory result due to Weyl for bounding the change in any singular value upon hermitian perturbation of a hermitian matrix. \blacksquare

Therefore, if we can approximate the matrix product $A^T A$, we immediately get a good reconstruction for every k . The appropriate sampling probabilities from the previous section are

$$p_t \geq \beta \frac{\mathbf{a}_t^T \mathbf{a}_t}{\|A\|_F^2}.$$

In this case, if $r \geq (4\rho/\beta\epsilon^2) \ln \frac{2d}{\delta}$, then with probability at least $1 - \delta$,

$$\|A - \tilde{A}_k\|^2 \leq \|A - A_k\|^2 + 2\epsilon\|A\|^2.$$

The sampling probabilities are easy to compute and sampling can be accomplished in one pass if the matrix is stored row-by-row.

To get a relative error result, we need a more carefully constructed set of non-uniform sampling probabilities. The problem here becomes apparent if A has rank k . In this case we have no hope of a relative error approximation unless we preserve the rank during sampling. To do so, we need to sample according to the actual singular vectors in U , not according to A ; this is because sampling according to A can give especially large weight to a few of the large singular value directions, ignoring the small singular value directions and hence not preserving rank. By sampling according to U , we essentially put equal weight on all singular directions. To approximate U well, we need sampling probabilities

$$p_t \geq \frac{\beta}{d} \mathbf{u}_t^T \mathbf{u}_t.$$

Then, from Corollary 13, if $r \geq (4(d - \beta)/\beta\epsilon^2) \ln \frac{2d}{\delta}$, with probability at least $1 - \delta$,

$$\|I - U^T Q^T Q U\| \leq \epsilon.$$

Since $\|U\| = 1$, it also follows that

$$\|UU^T - UU^T Q^T Q UU^T\| \leq \epsilon.$$

This result is useful because of the following lemma.

Lemma 20 (Spielman and Srivastava (2008)). *If $\|UU^T - UU^T Q^T Q UU^T\| \leq \epsilon$, then for every $\mathbf{x} \in \mathbb{R}^d$,*

$$(1 - \epsilon) \mathbf{x}^T A^T A \mathbf{x} \leq \mathbf{x}^T \tilde{A}^T \tilde{A} \mathbf{x} \leq (1 + \epsilon) \mathbf{x}^T A^T A \mathbf{x}.$$

Proof. We give a sketch of the proof from Spielman and Srivastava (2008). We let $\mathbf{x} \neq \mathbf{0}$ range over $\text{col}(U)$. Since $\text{col}(U) = \text{col}(A)$, $\mathbf{x} \in \text{col}(U)$ if and only if for some $\mathbf{y} \in \mathbb{R}^d$, $\mathbf{x} = A\mathbf{y}$. Since $\text{rank}(A) = d$, $A\mathbf{y} \neq 0 \iff \mathbf{y} \neq 0$. Also note that $UU^T A = A$, since UU^T is a projection operator onto the column space of U , which is the same as the column space of A . The following sequence establishes the lemma.

$$\begin{aligned} \|UU^T - UU^T Q^T Q UU^T\| &= \sup_{\mathbf{x} \neq \mathbf{0}} \frac{|\mathbf{x}^T UU^T \mathbf{x} - \mathbf{x}^T UU^T Q^T Q UU^T \mathbf{x}|}{\mathbf{x}^T \mathbf{x}}, \\ &= \sup_{A\mathbf{y} \neq \mathbf{0}} \frac{|\mathbf{y}^T A^T UU^T A \mathbf{y} - \mathbf{y}^T A^T UU^T Q^T Q UU^T A \mathbf{y}|}{\mathbf{y}^T A^T A \mathbf{y}}, \\ &= \sup_{A\mathbf{y} \neq \mathbf{0}} \frac{|\mathbf{y}^T A^T A \mathbf{y} - \mathbf{y}^T A^T Q^T Q A \mathbf{y}|}{\mathbf{y}^T A^T A \mathbf{y}}, \\ &= \sup_{\mathbf{y} \neq \mathbf{0}} \frac{|\mathbf{y}^T A^T A \mathbf{y} - \mathbf{y}^T \tilde{A}^T \tilde{A} \mathbf{y}|}{\mathbf{y}^T A^T A \mathbf{y}}, \end{aligned}$$

The lemma now follows because $\|UU^T - UU^T Q^T Q UU^T\| \leq \epsilon$. ■

Via the Courant-Fischer characterization Golub and Van Loan (1983) of the singular values, it is immediate from Lemma 20 that the singular value spectrum is also preserved :

$$(1 - \epsilon)\sigma_i(A^T A) \leq \sigma_i(\tilde{A}^T \tilde{A}) \leq (1 + \epsilon)\sigma_i(A^T A). \quad (5)$$

Lemma 20 along with (5) will allow us to prove the relative approximation result.

Theorem 21. *If $p_t \geq \frac{\beta}{d} \mathbf{u}_t^T \mathbf{u}_t$ and $r \geq (4(d - \beta)/\beta\epsilon^2) \ln \frac{2d}{\epsilon}$, then, for $k = 1, \dots, d$,*

$$\|A - A\tilde{\Pi}_k\| \leq \left(\frac{1 + \epsilon}{1 - \epsilon}\right)^{1/2} \|A - A_k\|,$$

where $\tilde{\Pi}_k$ projects onto the top k right singular vectors of \tilde{A} .

Remarks For $\epsilon \leq \frac{1}{2}$, $\left(\frac{1 + \epsilon}{1 - \epsilon}\right)^{1/2} \leq 1 + 2\epsilon$. Computing the probabilities p_t involves knowing \mathbf{u}_t which means one has to perform an *SVD*, in which case, one could use A_k ; it seems like overkill to compute A_k in order to approximate A_k . We discuss approximate sampling schemes later, in Section 7.

Proof. Let $\|\mathbf{x}\| = 1$. The following sequence establishes the result.

$$\begin{aligned} \|A(I - \tilde{\Pi}_k)\|^2 &= \sup_{\mathbf{x} \in \ker(\tilde{\Pi}_k)} \|A\mathbf{x}\|^2 = \sup_{\mathbf{x} \in \ker(\tilde{\Pi}_k)} \mathbf{x}^T A^T A \mathbf{x}, \\ &\leq \frac{1}{1 - \epsilon} \sup_{\mathbf{x} \in \ker(\tilde{\Pi}_k)} \mathbf{x}^T \tilde{A}^T \tilde{A} \mathbf{x}, \\ &= \frac{1}{1 - \epsilon} \sigma_{k+1}(\tilde{A}^T \tilde{A}), \\ &\leq \frac{1 + \epsilon}{1 - \epsilon} \sigma_{k+1}(A^T A) = \frac{1 + \epsilon}{1 - \epsilon} \|A - A_k\|^2. \end{aligned}$$

■

6 ℓ_2 Linear Regression with Relative Error Bounds

A linear regression is represented by a real data matrix $A \in \mathbb{R}^{m \times d}$ which represents m points in \mathbb{R}^d , and a target vector $\mathbf{y} \in \mathbb{R}^m$. Traditionally, $m \gg d$ (severly over constrained regression). The goal is to find a regression vector $\mathbf{x}^* \in \mathbb{R}^d$ which minimizes the ℓ_2 fit error (least squares regression)

$$\mathcal{E}(\mathbf{x}) = \|A\mathbf{x} - \mathbf{y}\|_2^2 = \sum_{t=1}^m (\mathbf{a}_t^T \mathbf{x} - y_t)^2,$$

We assume such an optimal \mathbf{x}^* exists (it may not be unique unless A has full column rank), and is given by $\mathbf{x}^* = A^+ \mathbf{y}$, where $^+$ denotes the More-Penrose pseudo-inverse; this problem can be solved in $O(md^2)$. Through row-sampling, it is possible to construct $\hat{\mathbf{x}}$, an approximation to the optimal regression weights \mathbf{x}^* , which is a relative error approximation to optimal,

$$\mathcal{E}(\hat{\mathbf{x}}) \leq (1 + \epsilon)\mathcal{E}(\mathbf{x}^*).$$

As usual, let $A = U_A S_A V_A^T$. Then $A^+ = V_A S_A^{-1} U_A^T$, and so $\mathbf{x}^* = V S^{-1} U^T \mathbf{y}$. The predictions are $\mathbf{y}^* = A \mathbf{x}^* = U_A U_A^T \mathbf{y}$, which is the projection of \mathbf{y} onto the column space of A . We define the residual $\boldsymbol{\epsilon} = \mathbf{y} - \mathbf{y}^* = \mathbf{y} - A \mathbf{x}^* = (\mathbf{I} - U_A U_A^T) \mathbf{y}$, so

$$\mathbf{y} = U_A U_A^T \mathbf{y} + \boldsymbol{\epsilon}. \quad (6)$$

We will construct \tilde{A} and $\tilde{\mathbf{y}}$ by sampling rows:

$$[\tilde{A}, \tilde{\mathbf{y}}] = Q[A, \mathbf{y}],$$

and solve the linear regression problem on $(\tilde{A}, \tilde{\mathbf{y}})$ to obtain $\hat{\mathbf{x}} = \tilde{A}^+ \tilde{\mathbf{y}}$. For $\beta \in (0, \frac{1}{3}]$, we will use the sampling probabilities

$$p_t \geq \beta \left(\frac{\mathbf{u}_t^2}{d} + \frac{(\mathbf{u}_t^2 + \frac{\epsilon_t^2}{\boldsymbol{\epsilon}^T \boldsymbol{\epsilon}})}{d+1} + \frac{\epsilon_t^2}{\boldsymbol{\epsilon}^T \boldsymbol{\epsilon}} \right) \quad (7)$$

to construct \tilde{A} and $\tilde{\mathbf{y}}$. There are three parts to these sampling probabilities. The first part allows us to reconstruct A well from \tilde{A} ; the second allows us to reconstruct $A^T \boldsymbol{\epsilon}$; and, the third allows us to reconstruct $\boldsymbol{\epsilon}$.

Note that $\tilde{A} = Q U_A S_A V_A^T$; if $Q U_A$ consisted of orthonormal columns, then this would be the SVD of \tilde{A} . Indeed, this is approximately so, as we will soon see. Let the SVD of \tilde{A} be $\tilde{A} = U_{\tilde{A}} S_{\tilde{A}} V_{\tilde{A}}^T$. Let $\tilde{U} = Q U_A$. Since $p_t \geq \beta \mathbf{u}_t^2/d$, it follows from Corollary 13 that if $r \geq 2 \frac{d-\beta}{\beta \epsilon^2}$, for $\epsilon \in (0, 1)$, then, with high probability,

$$\|\mathbf{I} - \tilde{U}^T \tilde{U}\| \leq \epsilon.$$

Since the eigenvalues of $\mathbf{I} - \tilde{U}^T \tilde{U}$ are given by $1 - \sigma_i^2(\tilde{U})$, it follows that

$$1 - \epsilon < \sigma_i^2(\tilde{U}) < 1 + \epsilon.$$

So all the singular values of U_A are preserved after sampling. Essentially, it suffices to sample $r = O(d \ln d / \epsilon^2)$ rows to preserve the entire spectrum of U_A . By choosing (say) $\epsilon = \frac{1}{2}$, the rank of U_A is preserved with high probability, since all the singular values are bigger than $\frac{1}{2}$. Thus, with high probability, $\text{rank}(\tilde{A}) = \text{rank}(U_{\tilde{A}}) = \text{rank}(Q U_A) = \text{rank}(U_A) = \text{rank}(A)$. Since $Q U_A$ has full rank, $S_{Q U_A}^{-1}$ is defined, and $S_{Q U_A} - S_{Q U_A}^{-1}$ is a diagonal matrix whose diagonals are $(\sigma_i^2(\tilde{U}) - 1)/\sigma_i(\tilde{U})$; thus, $\|S_{Q U_A} - S_{Q U_A}^{-1}\|_2 \leq \epsilon/\sqrt{1-\epsilon}$. This allows us to quantify the degree to which $Q U_A$ is orthonormal, because

$$\begin{aligned} \|(Q U_A)^+ - (Q U_A)^T\|_2 &= \|V_{Q U_A} S_{Q U_A}^{-1} U_{Q U_A}^T - V_{Q U_A} S_{Q U_A} U_{Q U_A}^T\|_2 \\ &= \|S_{Q U_A}^{-1} - S_{Q U_A}\|_2 \leq \frac{\epsilon}{\sqrt{1-\epsilon}}. \end{aligned}$$

Finally, we can get a convenient form for $\tilde{A}^+ = (Q A)^+$, because $Q A = Q U_A S_A V_A^T$ has full rank, and so $Q U_A = U_{Q U_A} S_{Q U_A} V_{Q U_A}^T$ has full rank (and hence is the product of full rank matrices). Thus,

$$\begin{aligned} (Q A)^+ &= (U_{Q U_A} S_{Q U_A} V_{Q U_A}^T S_A V_A^T)^+, \\ &= V_A (S_{Q U_A} V_{Q U_A}^T S_A)^+ U_{Q U_A}^T, \\ &= V_A S_A^{-1} V_{Q U_A} S_{Q U_A}^{-1} U_{Q U_A}^T, \\ &= V_A S_A^{-1} (Q U_A)^+, \end{aligned}$$

We summarize all this information in the next lemma.

Lemma 22. *If $r \geq (4d/\beta\epsilon^2) \ln \frac{2d}{\delta}$, with probability at least $1 - \delta$, all of the following hold:*

$$\text{rank}(\tilde{A}) = \text{rank}(U_{\tilde{A}}) = \text{rank}(QU_A) = \text{rank}(U_A) = \text{rank}(A); \quad (8)$$

$$\|S_{QU_A} - S_{QU_A}^{-1}\|_2 \leq \epsilon/\sqrt{1-\epsilon}; \quad (9)$$

$$\|(QU_A)^+ - (QU_A)^T\|_2 \leq \epsilon/\sqrt{1-\epsilon}; \quad (10)$$

$$(QA)^+ = V_A S_A^{-1} (QU_A)^+. \quad (11)$$

In Lemma 22 we have simplified the constant to 4; this is a strengthened form of Lemma 4.1 in Drineas *et al.* (2006d); in particular, the dependence on d is near-linear.

Remember that $\hat{\mathbf{x}} = \tilde{A}^+ \tilde{\mathbf{y}}$; we now bound $\|A\hat{\mathbf{x}} - \mathbf{y}\|^2$. We only sketch the derivation which basically follows the line of reasoning in Drineas *et al.* (2006d). Under the conditions of Lemma 22, with probability at least $1 - \delta$,

$$\begin{aligned} \|A\hat{\mathbf{x}} - \mathbf{y}\| &= \|A\tilde{A}^+ \tilde{\mathbf{y}} - \mathbf{y}\| = \|A(QA)^+ Q\mathbf{y} - \mathbf{y}\| \\ &\stackrel{(a)}{=} \|U_A(QU_A)^+ Q\mathbf{y} - \mathbf{y}\| \\ &\stackrel{(b)}{=} \|U_A(QU_A)^+ Q(U_A U_A^T \mathbf{y} + \boldsymbol{\epsilon}) - U_A U_A^T \mathbf{y} - \boldsymbol{\epsilon}\| \\ &\stackrel{(c)}{=} \|U_A(QU_A)^+ Q\boldsymbol{\epsilon} - \boldsymbol{\epsilon}\| \\ &= \|U_A((QU_A)^+ - (QU_A)^T)Q\boldsymbol{\epsilon} + U_A(QU_A)^T Q\boldsymbol{\epsilon} - \boldsymbol{\epsilon}\| \\ &\stackrel{(d)}{\leq} \|((QU_A)^+ - (QU_A)^T)\| \|Q\boldsymbol{\epsilon}\| + \|U_A^T Q^T Q\boldsymbol{\epsilon}\| + \|\boldsymbol{\epsilon}\| \\ &\stackrel{(e)}{\leq} \frac{\epsilon}{\sqrt{1-\epsilon}} \|Q\boldsymbol{\epsilon}\| + \|U_A^T Q^T Q\boldsymbol{\epsilon}\| + \|\boldsymbol{\epsilon}\|. \end{aligned}$$

(a) follows from Lemma 22; (b) follows from (6); (c) follows Lemma 22, because QU_A has full rank and so $(QU_A)^+ QU_A = I_d$; (d) follows from the triangle inequality and sub-multiplicativity using $\|U_A\| = 1$; finally, (e) follows from Lemma 22. We now see the rationale for the complicated sampling probabilities. Since $p_t \geq \epsilon_t^2/\epsilon^T \boldsymbol{\epsilon}$, for r large enough, by Theorem 17, $\|Q\boldsymbol{\epsilon}\|^2 \leq \|\boldsymbol{\epsilon}\|^2(1+\epsilon)$. Similarly, since $U_A^T \boldsymbol{\epsilon} = 0$, $\|U_A^T Q^T Q\boldsymbol{\epsilon}\| = \|U_A^T \boldsymbol{\epsilon} - U_A^T Q^T Q\boldsymbol{\epsilon}\|$; so, we can apply Lemma 14 with $S_1 = I_d$, $V = \boldsymbol{\epsilon}/\|\boldsymbol{\epsilon}\|$ and $S_2 = \|\boldsymbol{\epsilon}\|$. According to Lemma 14, if $p_t \geq \beta(\mathbf{u}_t^2 + \epsilon_t^2/\epsilon^T \boldsymbol{\epsilon})/(d+1)$, then if r is large enough, $\|U_A^T Q^T Q\boldsymbol{\epsilon}\| \leq \epsilon\|\boldsymbol{\epsilon}\|$. Since these are all probabilistic statements, we need to apply the union bound to ensure that all of them hold. Ultimately, we have the claimed result:

Theorem 23. *For sampling probabilities satisfying (7), and for $r \geq (8(d+1)/\beta\epsilon^2) \ln \frac{2(d+1)}{\delta}$, let $\hat{\mathbf{x}} = (QA)^+ Q\mathbf{y}$ be the approximate regression. Then, with probability at least $1 - 3\delta$,*

$$\|A\hat{\mathbf{x}} - \mathbf{y}\| \leq \left(1 + \epsilon + \epsilon\sqrt{\frac{1+\epsilon}{1-\epsilon}}\right) \|A\mathbf{x}^* - \mathbf{y}\|,$$

where $\mathbf{x}^* = A^+ \mathbf{y}$ is the optimal regression.

Remarks For the proof of the theorem, we observe that any transformation matrix Q satisfying the following three properties with high probability will do:

$$(i) \|I - U^T Q^T Q U\| \leq \epsilon; \quad (ii) \|Q\boldsymbol{\epsilon}\| \leq (1+\epsilon)\|\boldsymbol{\epsilon}\|; \quad (iii) \|U^T Q^T Q\boldsymbol{\epsilon}\| \leq \epsilon\|\boldsymbol{\epsilon}\|.$$

7 Estimating the Spectral Norm

The row-norm based sampling is relatively straightforward for the symmetric product. For the asymmetric product, $A^T B$, we need probabilities

$$p_t \geq \beta \frac{\frac{1}{\|A\|^2} \mathbf{a}_t^T \mathbf{a}_t + \frac{1}{\|B\|^2} \mathbf{b}_t^T \mathbf{b}_t}{\rho_A + \rho_B}. \quad (12)$$

To get these probabilities, we need $\|A\|$ and $\|B\|$; since we can compute the exact product in $O(md_1 d_2)$, a practically useful algorithm would need to estimate $\|A\|$ and $\|B\|$ efficiently. Suppose we had estimates λ_A, λ_B which satisfy:

$$(1 - \epsilon)\|A\|^2 \leq \lambda_A^2 \leq (1 + \epsilon)\|A\|^2; \quad (1 - \epsilon)\|B\|^2 \leq \lambda_B^2 \leq (1 + \epsilon)\|B\|^2.$$

We can construct probabilities satisfying the desired property with $\beta = (1 - \epsilon)/(1 + \epsilon)$.

$$\begin{aligned} p_t &= \frac{\frac{1}{\lambda_A^2} \mathbf{a}_t^T \mathbf{a}_t + \frac{1}{\lambda_B^2} \mathbf{b}_t^T \mathbf{b}_t}{\|A\|_F^2 / \lambda_A^2 + \|B\|_F^2 / \lambda_B^2} \\ &\geq \frac{\frac{1}{(1+\epsilon)\|A\|^2} \mathbf{a}_t^T \mathbf{a}_t + \frac{1}{(1+\epsilon)\|A\|^2} \mathbf{b}_t^T \mathbf{b}_t}{\|A\|_F^2 / (1 - \epsilon)\|A\|^2 + \|B\|_F^2 / (1 - \epsilon)\|A\|^2} \\ &= \left(\frac{1 - \epsilon}{1 + \epsilon} \right) \frac{\frac{1}{\|A\|^2} \mathbf{a}_t^T \mathbf{a}_t + \frac{1}{\|B\|^2} \mathbf{b}_t^T \mathbf{b}_t}{\rho_A + \rho_B}. \end{aligned}$$

One practical way to obtain $\|A\|^2$ is using the power iteration. Given an arbitrary unit vector \mathbf{x}_0 , for $n \geq 1$, let $\mathbf{x}_n = A^T A \mathbf{x}_{n-1} / \|A^T A \mathbf{x}_{n-1}\|$. Note that multiplying by $A^T A$ can be done in $O(2md_1)$ operations. Since \mathbf{x}_n is a unit vector, $\|A^T A \mathbf{x}_n\| \leq \|A\|^2$. We now get a lower bound. Let \mathbf{x}_0 be a random isotropic vector constructed using d_1 independent standard Normal variates z_1, \dots, z_{d_1} ; so $\mathbf{x}_0^T = [z_1, \dots, z_{d_1}] / \sqrt{z_1^2 + \dots + z_{d_1}^2}$. Let $\lambda_n^2 = \|A^T A \mathbf{x}_n\|$ be an estimate for $\|A\|^2$ after n power iterations.

Lemma 24. *For some constant $c \leq (\frac{2}{\pi} + 2)^3$, with probability at least $1 - \delta$,*

$$\lambda_n^2 \geq \frac{\|A\|^2}{\sqrt{4 + \frac{cd_1}{\delta^3}} \cdot 2^{-2n}}.$$

Remarks $n \geq c \log \frac{d_1}{\delta}$ gives the desired constant factor approximation. Since each power iteration takes $O(md_1)$ time, and we run $O(\log \frac{d_1}{\delta})$ power iterations, in $O(md_1 \log \frac{d_1}{\delta})$ time, we obtain a sufficiently good estimate for $\|A\|$ (and similarly for $\|B\|$).

Proof. Assume that $\mathbf{x}_0 = \sum_{i=1}^{d_1} \alpha_i \mathbf{v}_i$, where \mathbf{v}_i are the eigenvectors of $A^T A$ with corresponding eigenvalues $\sigma_1^2 \geq \dots \geq \sigma_{d_1}^2$. Note $\|A\|^2 = \sigma_1^2$. If $\sigma_{d_1}^2 \geq \sigma_1^2/2$, then it trivially follows that $\|A^T A \mathbf{x}_n\| \geq \sigma_1^2/2$ for any n , so assume that $\sigma_{d_1}^2 < \sigma_1^2/2$. We can thus partition the singular values into those at least $\sigma_1^2/2$ and those which are smaller; the latter set is non-empty. So assume for some $k < d_1$, $\sigma_k^2 \geq \sigma_1^2/2$ and $\sigma_{k+1}^2 < \sigma_1^2/2$. Since $\mathbf{x}_n = \sum_i \alpha_i \sigma_i^{2n} \mathbf{v}_i / (\sum_i \alpha_i^2 \sigma_i^{4n})^{1/2}$, we therefore

have:

$$\begin{aligned}
\lambda_n^4 &= \|\mathbf{A}^\top \mathbf{A} \mathbf{x}_n\|^2 = \frac{\sum_{i=1}^{d_1} \alpha_i^2 \sigma_i^{4(n+1)}}{\sum_{i=1}^{d_1} \alpha_i^2 \sigma_i^{4n}} \\
&\geq \frac{\sum_{i=1}^k \alpha_i^2 \sigma_i^{4(n+1)}}{\sum_{i=1}^{d_1} \alpha_i^2 \sigma_i^{4n}} \\
&= \frac{\sum_{i=1}^k \alpha_i^2 \sigma_i^{4(n+1)}}{\sum_{i=1}^k \alpha_i^2 \sigma_i^{4n} + \sum_{i=k+1}^{d_1} \alpha_i^2 \sigma_i^{4n}}, \\
&= \sigma_1^4 \frac{\sum_{i=1}^k \alpha_i^2 (\sigma_i/\sigma_1)^{4(n+1)}}{\sum_{i=1}^k \alpha_i^2 (\sigma_i/\sigma_1)^{4n} + \sum_{i=k+1}^{d_1} \alpha_i^2 (\sigma_i/\sigma_1)^{4n}}, \\
&\stackrel{(a)}{\geq} \sigma_1^4 \frac{\sum_{i=1}^k \alpha_i^2 (\sigma_i/\sigma_1)^{4(n+1)}}{4 \sum_{i=1}^k \alpha_i^2 (\sigma_i/\sigma_1)^{4(n+1)} + 2^{-2n}}, \\
&= \frac{\sigma_1^4}{4 + 2^{-2n} / \sum_{i=1}^k \alpha_i^2 (\sigma_i/\sigma_1)^{4(n+1)}}, \\
&\stackrel{(b)}{\geq} \frac{\sigma_1^4}{4 + 2^{-2n} / \alpha_1^2}.
\end{aligned}$$

(a) follows because for $i \geq k+1$, $\sigma_i^2 < \sigma_1^2/2$; for $i \leq k$, $\sigma_1^2/\sigma_i^2 \leq 4$; and $\sum_{i \geq k+1} \alpha_i^2 \leq \sum_{i \geq 1} \alpha_i^2 = 1$.

(b) follows because $\sum_{i=1}^k \alpha_i^2 (\sigma_i/\sigma_1)^{4(n+1)} \geq \alpha_1^2$. The theorem will now follow if we show that with probability at least $1 - c\delta^{1/3}$, $\alpha_1^2 \geq \delta/d$. It is clear that $\mathbb{E}[\alpha_1^2] = 1/d$ from isotropy. Without loss of generality, assume \mathbf{v}_1 is aligned with the z_1 axis. So $\alpha_1^2 = z_1^2 / \sum_i z_i^2$ (z_1, \dots, z_d are independent standard normals). For $\delta < 1$, we estimate $\mathbb{P}[\alpha_1^2 \geq \delta/d]$ as follows:

$$\begin{aligned}
\mathbb{P}\left[\alpha_1^2 \geq \frac{\delta}{d}\right] &= \mathbb{P}\left[\frac{z_1^2}{\sum_i z_i^2} \geq \frac{\delta}{d}\right] = \mathbb{P}\left[z_1^2 \geq \frac{\delta}{d} \sum_i z_i^2\right] = \mathbb{P}\left[z_1^2 \geq \frac{\delta}{d-\delta} \sum_{i \geq 2} z_i^2\right] \\
&\geq \mathbb{P}\left[z_1^2 \geq \frac{\delta}{d-1} \sum_{i \geq 2} z_i^2\right] \\
&\stackrel{(a)}{=} \mathbb{P}\left[\chi_1^2 \geq \frac{\delta}{d-1} \chi_{d-1}^2\right], \\
&\stackrel{(b)}{\geq} \mathbb{P}\left[\chi_1^2 \geq \delta + \delta^{2/3}\right] \cdot \mathbb{P}\left[\frac{\delta}{d-1} \chi_{d-1}^2 \leq \delta + \delta^{2/3}\right].
\end{aligned}$$

In (a) we compute the probability that a χ_1^2 random variable exceeds a multiple of an independent χ_{d-1}^2 random variable, which follows from the definition of the χ^2 distribution as a sum of squares of independent standard normals. (b) follows from independence and because one particular realization of the event in (a) is when $\chi_1^2 \geq \delta + \delta^{2/3}$ and $\delta \chi_{d-1}^2 / (d-1) \leq \delta + \delta^{2/3}$. Since $\mathbb{E}[\chi_{d-1}^2 / (d-1)] = 1$, and $\text{Var}[\chi_{d-1}^2 / (d-1)] = 2/(d-1)$, by Chebyshev's inequality,

$$\mathbb{P}\left[\frac{\delta}{d-1} \chi_{d-1}^2 \leq \delta + \delta^{2/3}\right] \geq 1 - \frac{2\delta^{1/3}}{d-1}.$$

From the definition of the χ_1^2 distribution, we can bound $\mathbb{P}[\chi_1^2 \leq \delta + \delta^{2/3}]$,

$$\mathbb{P}[\chi_1^2 \leq \delta + \delta^{2/3}] = \frac{1}{2^{1/2} \Gamma(1/2)} \int_0^{\delta + \delta^{2/3}} du \, u^{-1/2} e^{-u/2} \leq \sqrt{\frac{2}{\pi}} (\delta + \delta^{2/3})^{1/2},$$

and so

$$\mathbb{P} \left[\alpha_1^2 \geq \frac{\delta}{d} \right] \geq \left(1 - \sqrt{\frac{2}{\pi}} (\delta + \delta^{2/3})^{1/2} \right) \cdot \left(1 - \frac{2\delta^{1/3}}{d-1} \right) \geq 1 - \left(\frac{2}{\pi} + 2 \right) \delta^{1/3}.$$

■

We now consider the sampling based approach to estimate the spectral norm. Pre-sample the rows of A using probabilities proportional to the row norms to construct \tilde{A} . We know that if $r \geq (4\rho_A/\beta\epsilon^2) \ln \frac{2d_1}{\delta}$, then

$$\|\tilde{A}^T \tilde{A} - A^T A\| \leq \epsilon \|A\|^2.$$

It follows that we have a ϵ -approximation to the spectral norm from

$$\begin{aligned} \|\tilde{A}^T \tilde{A}\| &= \|\tilde{A}^T \tilde{A} - A^T A + A^T A\| \leq (1 + \epsilon) \|A\|^2; \\ \|A^T A\| &= \|A^T A - \tilde{A}^T \tilde{A} + \tilde{A}^T \tilde{A}\| \leq \epsilon \|A\|^2 + \|\tilde{A}^T \tilde{A}\|. \end{aligned}$$

Thus, $(1 - \epsilon) \|A\|^2 \leq \|\tilde{A}^T \tilde{A}\| \leq (1 + \epsilon) \|A\|^2$. Along this route, one must first sample r rows, and then approximate the spectral norm of the resulting \tilde{A} . We may now combine with the power iteration on $\tilde{A}^T \tilde{A}$ to get a constant factor approximation efficiently (or we may compute exactly in $O(rd_1^2)$). Specifically, set $\epsilon = \frac{1}{2}$, in which case, with high probability, $\frac{1}{2} \|A\|^2 \leq \|\tilde{A}^T \tilde{A}\| \leq \frac{3}{2} \|A\|^2$. Now, choose the number of power iterations $n \geq n^*$, where $\frac{\epsilon d_1}{\delta^3} = 2^{n^*}$. In this case, after n power iterations, we have an estimate which is at least $\frac{1}{2\sqrt{5}} \|\tilde{A}^T \tilde{A}\|^2$ from Lemma 24, which proves Theorem 25.

Theorem 25. *With $r \geq (4\rho_A/\epsilon^2) \ln \frac{2d_1}{\delta}$, the spectral norm estimate $\tilde{\sigma}_1^2$ obtained after $c \ln \frac{d_1}{\delta}$ power iterations on $\tilde{A}^T \tilde{A}$ starting from an isotropic random vector satisfies*

$$\frac{1}{2\sqrt{5}} \|A\|^2 \leq \tilde{\sigma}_1^2 \leq \frac{3}{2} \|A\|^2.$$

Further, the estimate $\tilde{\sigma}_1^2$ can be computed in $O(md_1 + \rho_A d_1 / \epsilon^2 \ln^2(\frac{d_1}{\delta}))$.

As mentioned at the begining of this section, constant factor approximations to the spectral norms of the relevant matrices is enough to obtain probabilities satisfying (12) for some constant β .

References

- Achlioptas, D. and McSherry, F. (2005). On spectral learning of mixtures of distributions. In *Proc. 18th Conference on Learning Theory (COLT)*, pages 458–469.
- Achlioptas, D. and McSherry, F. (2007). Fast computation of low-rank approximations. *Journal of ACM*, **54**(2), Article 10.
- Achlioptas, D., Fiat, A., Karlin, A., and McSherry, F. (2001). Web search via hub synthesis. In *Proc. 42nd IEEE symposium on Foundations of Computer Science*, pages 500–509.
- Ahlsvede, R. and Winter, A. (2002). Strong converse for identification via quantum channels. *Information Theory, IEEE Transactions on*, **48**(3), 569–579.
- Ailon, N. and Chazelle, B. (2006). Approximate nearest neighbors and the fast Johnson-Lindenstrauss transform. In *Proc. 38th ACM Symposium on Theory of Computing*, pages 557–563.

- Azar, Y., Fiat, A., Karlin, A., McSherry, F., and Saia, J. (2001). Spectral analysis of data. In *Proc. 33rd ACM symposium on Theory of computing*, pages 619–626.
- Bernstein, S. (1924). On a modification of chebyshev’s inequality and of the error formula of laplace. *Ann. Sci. Inst. Sav. Ukraine*, **4**(5), 38–49.
- Berry, M. W., Dumais, S. T., and O’Brien, G. W. (1995). Using linear algebra for intelligent information retrieval. *SIAM Rev.*, **37**(4), 573–595.
- Chernoff, H. (1952). A measure of asymptotic efficiency for tests of a hypothesis based on the sum of observations. *Annals of Mathematical Statistics*, **23**(4), 493–507.
- Deshpande, A. and Vempala, S. (2006). Adaptive sampling and fast low-rank matrix approximation. In *Proc. 10th RANDOM*.
- Deshpande, A., Rademacher, L., Vempala, S., and Wang, G. (2006). Matrix approximation and projective clustering via volume sampling. In *Proc. 17th ACM-SIAM symposium on Discrete algorithms*, pages 1117–1126.
- Drineas, P., Kerenidis, I., and Raghavan, P. (2002). Competitive recommendation systems. In *Proc. 34th ACM symposium on Theory of computing*, pages 82–90.
- Drineas, P., Frieze, A., Kannan, R., Vempala, S., and Vinay, V. (2004). Clustering large graphs via the singular value decomposition. *Machine Learning*, **56**(1-3), 9–33.
- Drineas, P., Kannan, R., and Mahoney, M. W. (2006a). Fast Monte Carlo algorithms for matrices i: Approximating matrix multiplication. *SIAM Journal on Computing*, **36**(1), 132–157.
- Drineas, P., Kannan, R., and Mahoney, M. W. (2006b). Fast Monte Carlo algorithms for matrices ii: Computing a low rank approximation to a matrix. *SIAM Journal on Computing*, **36**(1), 158–183.
- Drineas, P., Kannan, R., and Mahoney, M. W. (2006c). Fast Monte Carlo algorithms for matrices iii: Computing a compressed approximate matrix decomposition. *SIAM Journal on Computing*, **36**(1), 184–206.
- Drineas, P., Mahoney, M. W., and Muthukrishnan, S. (2006d). Sampling algorithms for l2 regression and applications. In *Proc. 17th ACM-SIAM Symposium on Discrete Algorithms (SODA)*, pages 1127–1136.
- Drineas, P., Mahoney, M. W., and Muthukrishnan, S. (2006e). Subspace sampling and relative-error matrix approximation: Column-row based methods. In *Proc. 14th ESA*.
- Drineas, P., Magdon-Ismail, M., and Mahoney, M. (2010). Sampling leverage scores.
- Frieze, A., Kannan, R., and Vempala, S. (1998). Fast monte-carlo algorithms for finding low-rank approximations. In *Proc. 39th Annual Symposium on Foundations of Computer Science*, pages 370–378.
- Golub, G. and Van Loan, C. (1983). *Matrix Computations*. Johns Hopkins University Press, Baltimore.
- Golub, G. and van Loan, C. (1996). *Matrix computations*. The Johns Hopkins University Press, London, 3 edition.

- Gross, D., Liu, Y.-K., Steven Flammia, S. T., Becker, S., and Eisert, J. (2009). Quantum state tomography via compressed sensing. preprint available at: arXiv:0909.3304v2.
- Hoeffding, W. (1963). Probability inequalities for sums of bounded random variables. *Journal of the American Statistical Association*, **58**(301), 13–30.
- Kannan, R., Salmasian, H., and Vempala, S. (2008). The spectral method for general mixture models. *SIAM Journal of Computing*, **38**(3), 1141–1156.
- Kleinberg, J. (1999). Authoritative sources in a hyperlinked environment. *Journal of the ACM*, **46**(5), 604–632.
- Kuczyński, J. and Woźniakowski, H. (1989). Estimating the largest eigenvalue by the power and lanczos algorithms with a random start. Technical Report CUCS-026-92, Institute of Informatics, University of Warsaw.
- Magdon-Ismail, M. (2010). Row sampling for matrix algorithms via a non-commutative bernstein bound. *CoRR*, **abs/1008.0587**.
- Magen, A. and Zouzias, A. (2010). Low rank matrix-valued chernoff bounds and applications. *submitted*. <http://arxiv.org/abs/1005.2724>.
- McSherry, F. (2001). Spectral partitioning of random graphs. In *Proc 42nd IEEE symposium on Foundations of Computer Science*, pages 529–537.
- Papadimitriou, C. H., Tamaki, H., Raghavan, P., and Vempala, S. (2000). Latent semantic indexing: A probabilistic analysis. *Journal of Computer and System Sciences*, **61**(2), 217–235.
- Recht, B. (2009). A simpler approach to matrix completion. working paper.
- Rudelson, M. and Vershynin, R. (2007). Sampling from large matrices: An approach through geometric functional analysis. *Journal of the ACM*, **54**(4), 19.
- Sarlos, T. (2006). Improved approximation algorithms for large matrices via random projections. In *Proc. 47th IEEE Symposium on Foundations of Computer Science*, pages 143–152.
- Spielman, D. A. and Srivastava, N. (2008). Graph sparsification by effective resistances. In *STOC '08: Proc. 40th ACM Symposium on Theory of Computing*, pages 563–568.
- Wigderson, A. and Xiao, D. (2008). Derandomizing the ahlsvede-winter matrix-valued chernoff bound using pessimistic estimators, and applications. *Theory of Computing*, **4**(3), 53–76.
- Woolfe, F., Liberty, E., Rokhlin, V., , and Tygert, M. (2008). A fast randomized algorithm for the approximation of matrices. *Applied and Computational Harmonic Analysis*, **25**(3), 335–366.